

ML Lab-2

Week-2

Name: Chandan Chatragadda	SRN: PES2UG23CS141	Section: 5C
---------------------------	--------------------	-------------

Week: 2

Mushrooms.csv

```
PS D:\Auss-TUP-CC-BOP-20240015\Chandan-Profile\Docs\PES\SPL\ML\all\code\sklearn_implementation> python test.py --D EC_C_PES2UG23CS141_L40 --data mushrooms.csv
Running tests with PYTORCH framework
=====
target column: 'class' (last column)
Original dataset info:
Shape: (8124, 23)
Columns: ['cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor', 'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root', 'stalk-surface-above-ring', 'stalk-surface-below-ring', 'stalk-color-above-ring', 'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number', 'ring-type', 'spore-print-color', 'population', 'habitat', 'class']
First few rows:
cap-shape: ['x' 'b' 's' 'f' 'x'] -> [5 0 4 2 3]
cap-surface: ['s' 'y' 'f' 'g'] -> [2 3 0 1]
cap-color: ['n' 'y' 'w' 'g' 'e'] -> [4 9 8 3 2]
class: ['p' 'e'] -> [1 0]
Processed dataset shape: torch.Size([8124, 23])
Number of features: 22
Features: ['cap-shape', 'cap-surface', 'cap-color', 'bruises', 'odor', 'gill-attachment', 'gill-spacing', 'gill-size', 'gill-color', 'stalk-shape', 'stalk-root', 'stalk-surface-above-ring', 'stalk-surface-below-ring', 'stalk-color-above-ring', 'stalk-color-below-ring', 'veil-type', 'veil-color', 'ring-number', 'ring-type', 'spore-print-color', 'population', 'habitat']
Target: class
Framework: PYTORCH
Data type: <class 'torch.Tensor'>
=====
DECISION TREE CONSTRUCTION LOG
=====
Total samples: 8124
Training samples: 6499
Testing samples: 1625
Constructing decision tree using training data...
Decision tree construction completed using PYTORCH
=====
OVERALL PERFORMANCE METRICS
=====
Accuracy: 1.0000 (100.00%)
Precision (weighted): 1.0000
Recall (weighted): 1.0000
F1-Score (weighted): 1.0000
Precision (macro): 1.0000
Recall (macro): 1.0000
F1-Score (macro): 1.0000
=====
TREE COMPLEXITY METRICS
=====
Maximum Depth: 4
Total Nodes: 29
Leaf Nodes: 26
Internal Nodes: 3
```

▲ DECISION TREE STRUCTURE

```
=====
Root [odor] (gain: 0.9083)
├── = 0:
│   └── Class 0
├── = 1:
│   └── Class 1
├── = 2:
│   └── Class 1
├── = 3:
│   └── Class 0
├── = 4:
│   └── Class 1
├── = 5:
│   ├── [spore-print-color] (gain: 0.1469)
│   ├── = 0:
│   │   └── Class 0
│   ├── = 1:
│   │   └── Class 0
│   ├── = 2:
│   │   └── Class 0
│   ├── = 3:
│   │   └── Class 0
│   ├── = 4:
│   │   └── Class 0
│   ├── = 5:
│   │   └── Class 1
│   └── = 7:
│       ├── [habitat] (gain: 0.2217)
│       ├── = 0:
│       │   ├── [gill-size] (gain: 0.7642)
│       │   ├── = 0:
│       │   │   └── Class 0
│       │   └── = 1:
│       │       └── Class 1
│       ├── = 1:
│       │   └── Class 0
│       ├── = 2:
│       │   ├── [cap-color] (gain: 0.7300)
│       │   ├── = 1:
│       │   │   └── Class 0
│       │   ├── = 4:
│       │   │   └── Class 0
│       │   ├── = 8:
│       │   │   └── Class 1
│       │   └── = 9:
│       │       └── Class 1
│       ├── = 4:
│       │   └── Class 0
│       ├── = 6:
│       │   └── Class 0
│       └── = 8:
│           └── Class 0
├── = 6:
│   └── Class 1
├── = 7:
│   └── Class 1
└── = 8:
    └── Class 1
```

Nursery.csv

```
PS D:\Asus-TUF-CC-BKP-20240915\Chandan-Profile\Docs\PES\SEM_5\ML\all\code\sklearn_implementation> python test.py --ID EC_C_PES2UG23CS141_Lab3 --data Nursery.csv
Running tests with PYTORCH framework
=====
target column: 'class' (last column)
Original dataset info:
Shape: (12960, 9)
Columns: ['parents', 'has_nurs', 'form', 'children', 'housing', 'finance', 'social', 'health', 'class']

First few rows:

parents: ['usual' 'pretentious' 'great_pret'] -> [2 1 0]
has_nurs: ['proper' 'less_proper' 'improper' 'critical' 'very_crit'] -> [3 2 1 0 4]
form: ['complete' 'completed' 'incomplete' 'foster'] -> [0 1 3 2]
class: ['recommend' 'priority' 'not_recom' 'very_recom' 'spec_prior'] -> [2 1 0 4 3]

Processed dataset shape: torch.Size([12960, 9])
Number of features: 8
Features: ['parents', 'has_nurs', 'form', 'children', 'housing', 'finance', 'social', 'health']
Target: class
Framework: PYTORCH
Data type: <class 'torch.Tensor'>

=====
DECISION TREE CONSTRUCTION DEMO
=====
Total samples: 12960
Training samples: 10368
Testing samples: 2592

Constructing decision tree using training data...

🌱 Decision tree construction completed using PYTORCH!

📊 OVERALL PERFORMANCE METRICS
=====
Accuracy:          0.9867 (98.67%)
Precision (weighted): 0.9876
Recall (weighted):  0.9867
F1-Score (weighted): 0.9872
Precision (macro):  0.7604
Recall (macro):     0.7654
F1-Score (macro):   0.7628

🌱 TREE COMPLEXITY METRICS
=====
Maximum Depth:      7
Total Nodes:        952
Leaf Nodes:         680
Internal Nodes:     272
```



Tictactoe.csv:

```
PS D:\Asus-UF-CC-BKP-20240915\Chandan-Profile\Docs\SES\SEM5\VL\all\code\sklearn_implementation> python test.py --data tictactoe.csv
Running tests with PYTORCH framework
=====
target column: 'Class' (last column)
Original dataset info:
Shape: (958, 10)
Columns: ['top-left-square', 'top-middle-square', 'top-right-square', 'middle-left-square', 'middle-middle-square', 'middle-right-square', 'bottom-left-square', 'bottom-middle-square', 'bottom-right-square', 'Class']

First few rows:

top-left-square: ['x' 'o' 'b'] -> [2 1 0]
top-middle-square: ['x' 'o' 'b'] -> [2 1 0]
top-right-square: ['x' 'o' 'b'] -> [2 1 0]

Class: ['positive' 'negative'] -> [1 0]

Processed dataset shape: torch.Size([958, 10])
Number of features: 9
Features: ['top-left-square', 'top-middle-square', 'top-right-square', 'middle-left-square', 'middle-middle-square', 'middle-right-square', 'bottom-left-square', 'bottom-middle-square', 'bottom-right-square']
Target: Class
Framework: PYTORCH
Data type: <class 'torch.Tensor'>
=====

DECISION TREE CONSTRUCTION DEMO
=====
Total samples: 958
Training samples: 766
Testing samples: 192

Constructing decision tree using training data...

🌲 Decision tree construction completed using PYTORCH!

📊 OVERALL PERFORMANCE METRICS
=====
Accuracy:          0.8730 (87.30%)
Precision (weighted): 0.8741
Recall (weighted):  0.8730
F1-Score (weighted): 0.8734
Precision (macro):  0.8590
Recall (macro):     0.8638
F1-Score (macro):   0.8613

🌲 TREE COMPLEXITY METRICS
=====
Maximum Depth:      7
Total Nodes:         281
Leaf Nodes:          180
Internal Nodes:       101
```

SECRET//NOFORN//REL//JAN

a) Algorithm Performance

Qn. Which dataset achieved the highest accuracy and why?

Ans: Mushroom dataset. As it has a properly defined dataset.

Qn. How does dataset size affect performance?

Ans: Bigger datasets have more information, boundaries and improve accuracy.

Qn. What role does the number of features play?

Ans: More features can make the model learn complex patterns and give better predictions and perform better in the real world datasets.

b) Data Characteristics Impact

Qn. How does class imbalance affect tree construction?

Ans: Class imbalance can bias the tree towards majority classes and give partial results.

Qn. Which types of features (binary vs multi-valued) work better?

Ans: Binary features often create simpler and clearer splits. Multi-valued categorical features offer richer information but

increases tree complexity.

c) Practical Applications

Qn. For which real-world scenarios is each dataset type most relevant?

Ans:

- **Mushroom: Food safety and scientific studies.**
- **Tic-tac-toe: Game outcome prediction or training a game bot.**
- **Nursery: Social and family-based recommendation systems with multiple categories reflecting real-world complexity.**

Qn. What are the interpretability advantages for each domain?

Ans: Mushroom dataset's binary classification and features give direct outcome of whether it is safe or poisonous.

Tictactoe dataset's game states map intuitively to tree splits and helps in understanding.

The nursery dataset, with multi-class and more complex categorical factors, offers richer but more complex interpretability challenges.

d) How to Improve Performance for Each Dataset

- **Mushroom:** Reduce overfitting, add feature selection and prune unnecessary features.

- **Tic-tac-toe:** Tune tree depth to prevent overfitting.

- **Nursery:** Handle class imbalance and consider hierarchical tree structures due to multi-class target and complex features.